

Supercomputer architectures

An Introduction to HPC

3 Nov 2015 | Alexander Peyser | Simulation Lab Neuroscience
Jülich Supercomputing Centre
Institute for Advanced Simulation
Forschungszentrum Jülich



What is High Performance Computing?

Hardware architectures

Current Systems

Working with Supercomputers

Programming Paradigms

What is High Performance Computing?

HPC definition

HPC capabilities

What do we have?

How do we use them?

Fully scaled projects

What is a Supercomputer?

Wikipedia: A supercomputer is a computer with a high-level computational capacity compared to a general-purpose computer. Performance of a supercomputer is measured in floating point operations per second (FLOPS) instead of million instructions per second (MIPS). As of 2015, there are supercomputers which can perform up to quadrillions of FLOPS



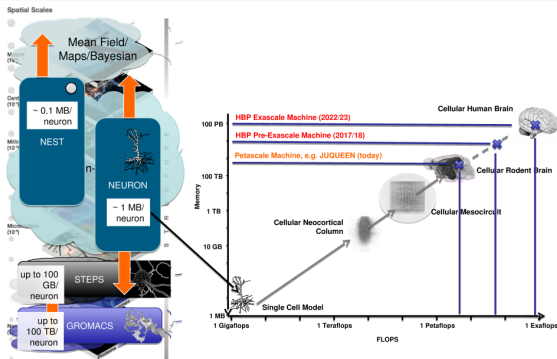
What is it in short?

A computer at the very limit of what is currently available



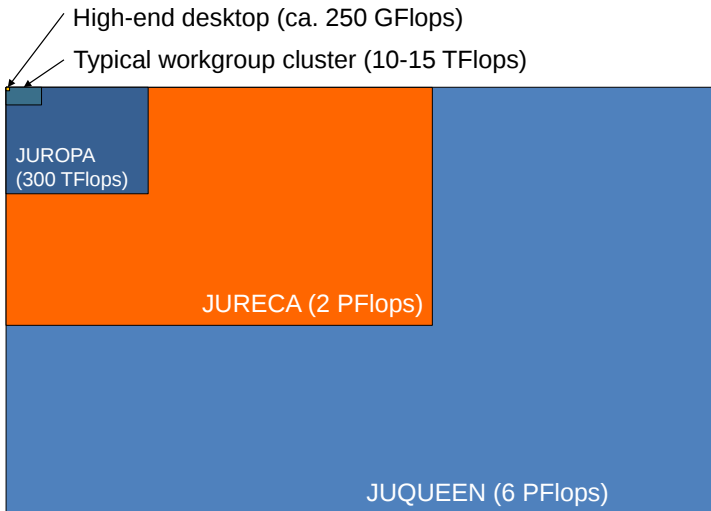
What do we do with Supercomputers?

Supercomputers are massively parallel machines using the most advanced nodes, interconnects and memory to tackle problems that are too large or take too long for commodity computers



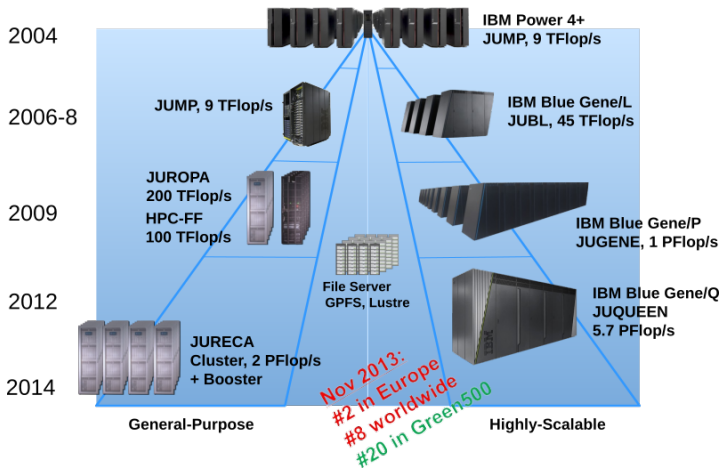
Felix Schürmann (EPFL)

Theoretical peak performance



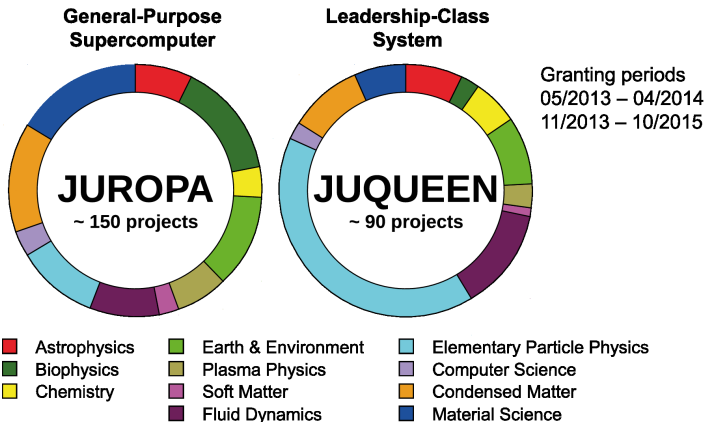
What do we have?

What is High Performance Computing?



How do we use them?

What is High Performance Computing?



Fully scaled projects

What is High Performance Computing?



CIAO
 Hydrolytic molecular dynamics solver for turbulent reacting flows in complex geometries.



CoreNeuron
 Scalable electrical activity of neuronal networks with morphologically detailed neurons.



dynQCD
 Lattice QCD on ChromoQuench (CQC) with dynamical fermions.



FE2T1 (ex. mFE²)
 A scale bridging approach incorporating micro-mechanics in macroscopic simulations of multi-phase solids.



FEQuash
 A framework for the massively parallel QM simulation of molecular systems governed by DFT.



Gysela
 A parallel code from CNRS, France, for modelling heterogeneous plasmas.



ICON
 A solver for fully compressible, nonhydrostatic equations of motion at very high-resolution resolution.



PMO-PPASST
 A multi-scale coupled solver for systems of QCD.



PP-Code
 A parallel multi-scale for simulating charged and neutral particle dynamics in astrophysical and non-relativistic plasmas.



proOpen
 Open numerical simulation of the solar heliosphere.



SHOCK
 Decoupled numerical simulation of astrophysical flows.



TerraNeo
 A multi-scale solver for geophysical applications from the University of Chicago.



velDoris
 A multi-scale multi-scale Lattice Boltzmann solver from the University of Stuttgart.



ZFS
 A multi-scale framework for compressible and incompressible flow, microphysics, and combustion phenomena.



JustPC
 A fully parallel Python-to-C++ code for wave equation simulation.



KORNano
 Kornelius-Bohrer (Korn) function code for quantum description of nanosystems.



LAMMPS (DCM)
 A Dynamic Load Method for a distributed molecular dynamics code, the Large-scale Atomic/Molecular Massively Parallel Simulator.



MPFC
 Fast simulations using a hybrid representation of reduced particles.



mp (muPh)
 Vector flow and scalar transport in porous media.



Musubi
 A multi-scale Lattice Boltzmann solver for flow simulation.



NEST
 A multi-scale multi-scale network model for the dynamics and evolution of neural systems.



JURASSIC
 A fast solver for turbulent cellular kinetics in the kinetic approach.



PEPC
 A parallel flow code for solving the Navier-Stokes equations for Compressible, unsteady and hydrodynamic systems.



IMD
 A software package for simulating molecular dynamics simulations.



OpenTBC
 Open parallel simulation of turbulent flows.

Hardware architectures

Von Neumann Architecture

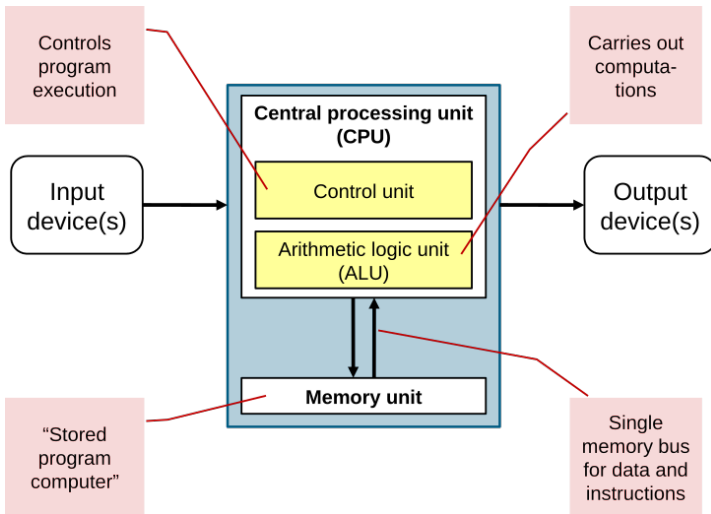
Von Neumann Bottlenecks

Parallelizing out of the bottlenecks

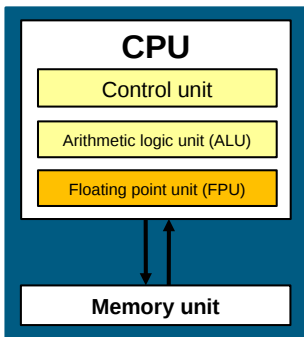
Parallelism speedup

General implementation

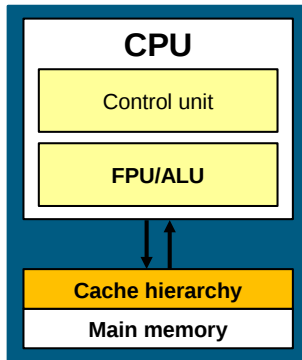
GPUs: data parallel problems



Computing with real numbers



Enhancing memory access



Flops floating point operations per second

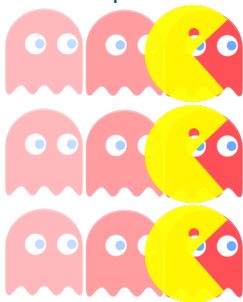
Bandwidth Number of bytes of memory transferred per second

Latency Time delay to completion of a single computational or memory access operation

One pacman eats nine ghosts in 3 seconds...



... but three pacmen eat 9 ghosts in 1 second



... which is strong scaling

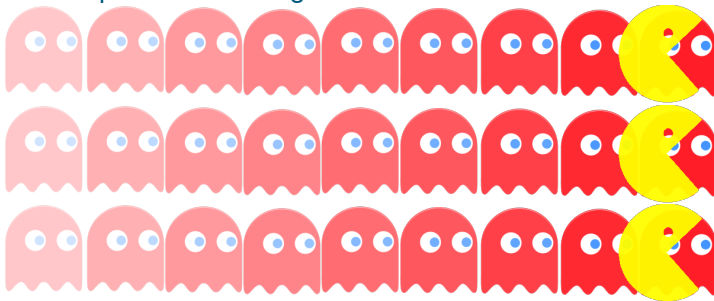
Parallelizing out of the bottlenecks

Hardware architectures

One pacman eats nine ghosts in 3 seconds...



... or three pacmen eat 27 ghosts in 3 seconds



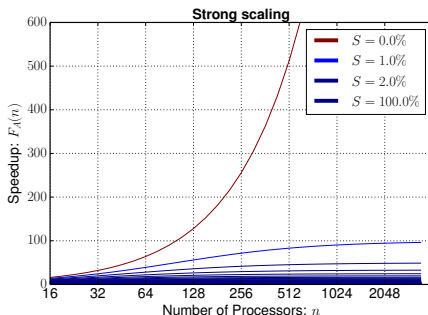
... which is weak scaling

Limits to parallelism

Amdahl's law: $F_A(n) = [S + (1 - S)/n]^{-1}$

If $S > 0$, $\lim_{n \rightarrow \infty} F_A(n) = 1/S$

Gustafson's law: $F_G(n) = S + n * (1 - S)$
 $(F_G(n) - S)/n = 1 - S$

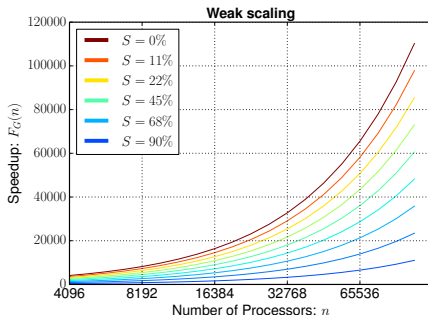


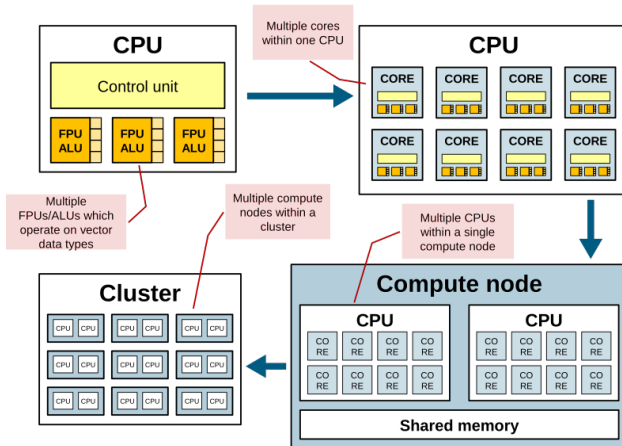
Limits to parallelism

Amdahl's law: $F_A(n) = [S + (1 - S)/n]^{-1}$

If $S > 0$, $\lim_{n \rightarrow \infty} F_A(n) = 1/S$

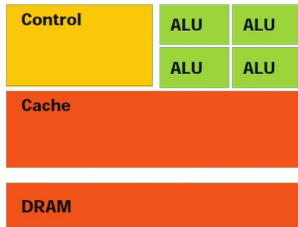
Gustafson's law: $F_G(n) = S + n * (1 - S)$
 $(F_G(n) - S)/n = 1 - S$





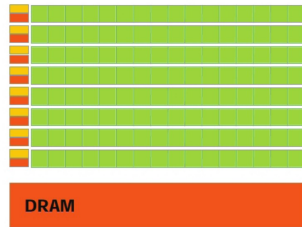
GPU: Graphics Processing Unit

A large number of arithmetic/floating point units with reduced control logic in parallel to the CPU. Originally developed for 3D graphics rendering



CPU

nvidia, amd



GPU

Current Systems

Performance

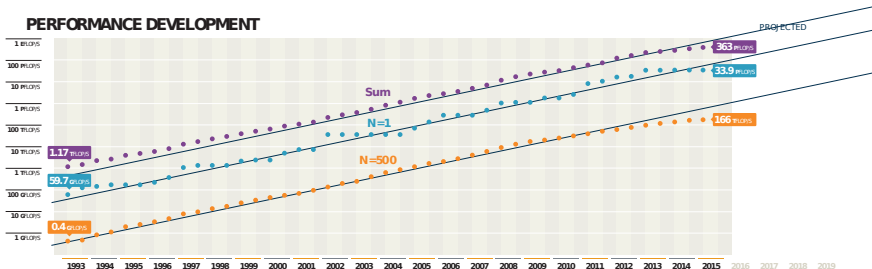
Architecture

Accelerators

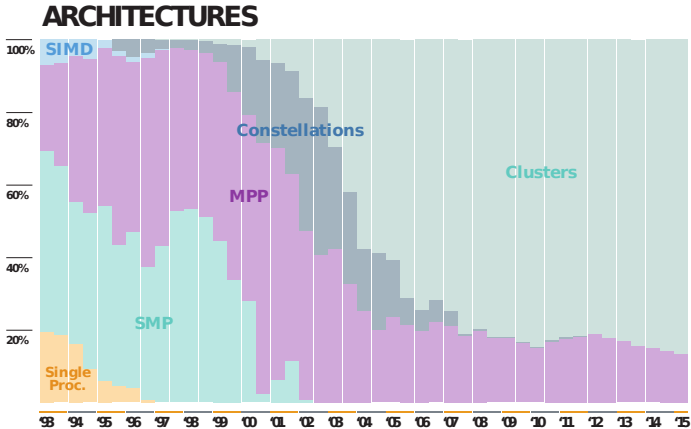
JUQUEEN

JURECA

PERFORMANCE DEVELOPMENT

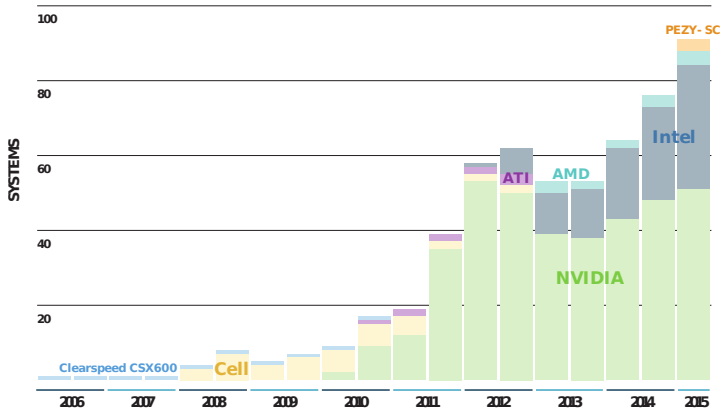


top500.org



top500.org

ACCELERATORS/CO-PROCESSORS



top500.org

Data Sheet	
Manufacturer	IBM
Design	Bluegene/Q
CPU Core	PowerPC A2 16 (x4)
Cores	458,752
Linpack Perf. (R_{max})	5.009 PFlop/s
Theoretical Peak (R_{peak})	5.872 PFlop/s
Power	2.3 MW
Memory	458 TB
Interconnect	Custom Interconnect Network
Topology	5D Torus
Operating System	Linux / CNK
Batch	LoadLeveler





Data Sheet

Manufacturer	T-Platforms, ParTec, Intel, Mellanox
Design	Intel / GPU
CPU Core	2 Intel Haswell 12-core & K80 GPUs
Cores	45,216
GPUs	75 nodes w/ 2 K80 (4992 CUDA)
Linpack Perf. (R_{max})	(Prel) 1.4 PFlop/s
Theoretical Peak (R_{peak})	1.8 PFlop/s + 0.44 (GPU)
Power	???
Memory	281 TB
Interconnect	Mellanox EDR InfiniBand
Operating System	CentOS 7 Linux
Batch	Slurm

Working with Supercomputers

Hardware

Cluster organization

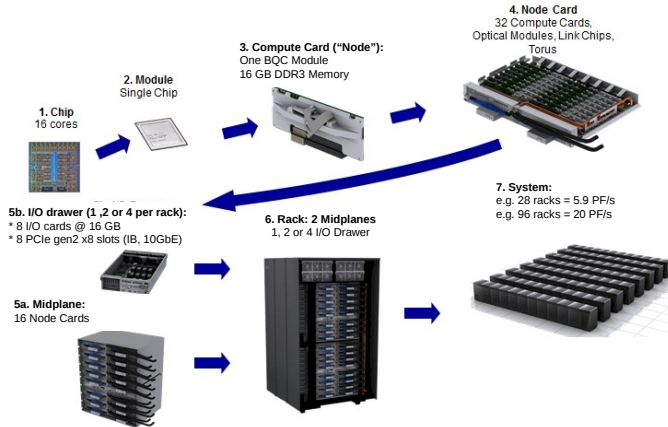
Running system

LoadLeveler commands

LoadLeveler script

Hardware

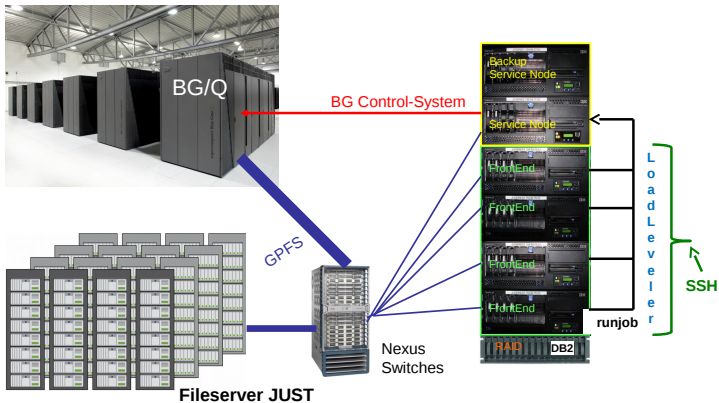
Working with Supercomputers



IBM

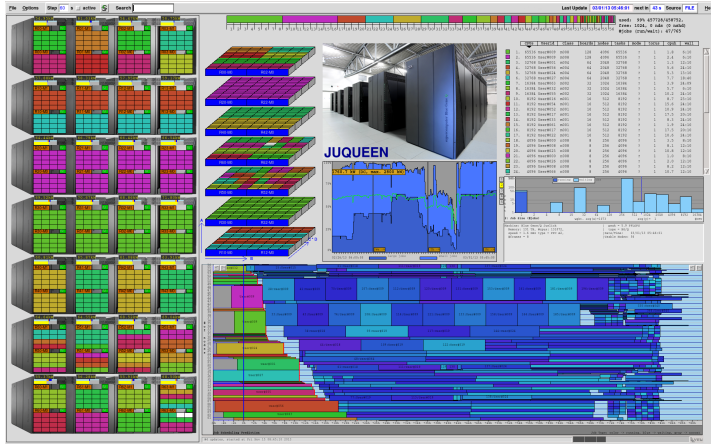
Cluster organization

Working with Supercomputers



Michael Stephan & Jutta Doctor (JSC)

Supercomputers are shared resources



Command	Description
<code>llsubmit <jobfile></code>	Send job to the queuing system
<code>llq</code>	List all queued and running jobs
<code>llq -l <job ID></code>	detail information about the specific job
<code>llq -s <job ID></code>	details information about a specific queued job, such as start time
<code>llq -u <user></code>	list all jobs from one user
<code>llcancel <job ID></code>	Kill the specified job
<code>llstatus</code>	Display the status of LoadLeveler
<code>llclass</code>	List existing classes and their properties
<code>llqx</code>	Show detailed information about all jobs

LoadLeveler script

Working with Supercomputers

```
# @job_name           = weakScale_nmr00032_nthr04_N16667_dryRun
# @job_type           = bluegene
# @bg_size            = 32
# @bg_connectivity    = TORUS
# @environment        = COPY_ALL
# @wall_clock_limit   = 00:30:00
# @output             = $(HOME)/log/juqueen/${job_name}.${jobid}.out
# @error              = $(HOME)/log/juqueen/${job_name}.${jobid}.err
# @notification       = error
# @notify_user        = w.schenck@fz-juelich.de
# @queue
```

```
export OMP_NUM_THREADS=4
```

```
NPROCS=128
NEST_VERSION=10kcollaps_dry_run
NEST_EXE=${HOME}/opt/NEST/juqueen_${NEST_VERSION}/bin/nest
SLI_SCRIPT=spike_counting_4G_dryRun.sli
```

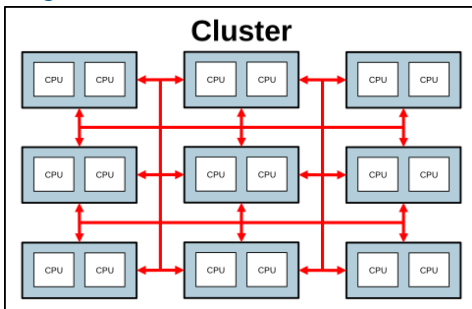
```
runjob --exe $NEST_EXE \\  
      --np $NPROCS \\  
      --ranks-per-node 1 \\  
      --verbose 1 --exp-env OMP_NUM_THREADS \\  
      --args $SLI_SCRIPT
```

Programming Paradigms

Overview

Languages

Tools

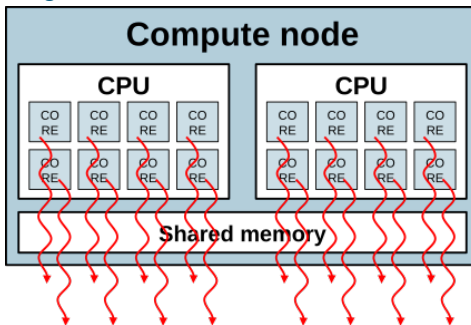


Decomposition

MPI At the cluster level, message passing standard via network between nodes (application)

OpenMP At the node level, thread programming standard using shared memory (process)

OpenCL At the GPU level, work groups of parallel kernels



Decomposition

MPI At the cluster level, message passing standard via network between nodes (application)

OpenMP At the node level, thread programming standard using shared memory (process)

OpenCL At the GPU level, work groups of parallel kernels

Overview

Programming Paradigms



Decomposition

MPI At the cluster level, message passing standard via network between nodes (application)

OpenMP At the node level, thread programming standard using shared memory (process)

OpenCL At the GPU level, work groups of parallel kernels

Which programming languages can I use?

Scripting Run over parameter spaces with scripting languages and JUBE — and run (most) any code

MPI C/C++/Fortran/Python/JAVA... (special library)

OpenMP C/C++/Fortran (special directives)

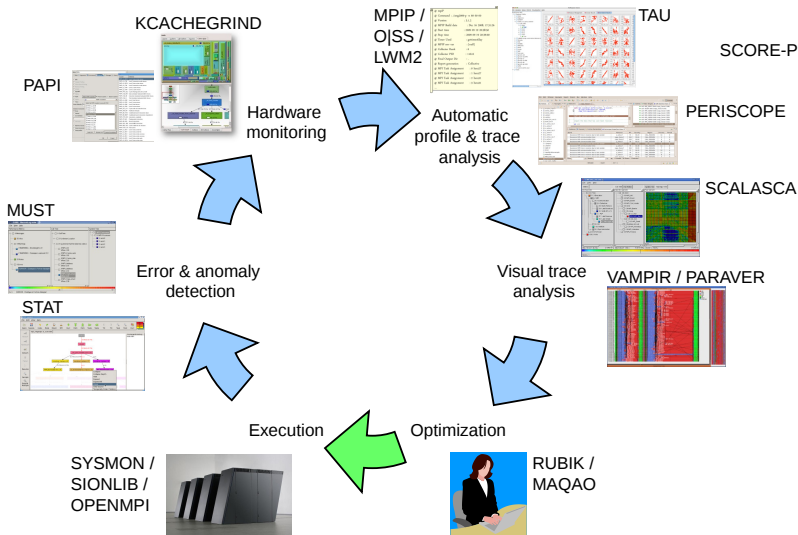
pthread Native threads for most languages

GPU C/C++ with other experimental kernel mappings, bindings on cpus to most everything... (kernels compiled and shipped)

Make the outer loops & code as you want, and focus on optimizing core kernels

Tools

Programming Paradigms



Introduction to Parallel Performance Engineering, VI-HPS

Overview Languages Tools

HPC Arch Systems Practical Programming Conclusions

Thank You!